

Explainable Deep Learning for Diabetic Retinopathy Detection Using Fundus Images

Lujain To'ma

Dept. Artificial intelligence

Jordan Univ. of Science and Technology

Irbid, Jordan

lytoma22@cit.just.edu.jo

Malak A. Abdullah

Dept. Computer Science

Jordan Univ. of Science and Technology

Irbid, Jordan

mabdullah@just.edu.jo

Areen AlAkaleek

Dept. Artificial intelligence

Jordan Univ. of Science and Technology

Irbid, Jordan

aaalakaleek22@cit.just.edu.jo

Toqa Aldagmseh

Dept. Artificial intelligence

Jordan Univ. of Science and Technology

Irbid, Jordan

tmaldagmsih22@cit.just.edu.jo

Abstract—Diabetic Retinopathy (DR) [18] is a leading cause of blindness worldwide, and early detection plays a critical role in preventing vision loss. In this research paper, we propose an explainable deep learning approach for DR detection using fundus images. Our approach leverages modern convolutional neural networks such as DenseNet121, known for their high accuracy in image classification tasks. To enhance the interpretability of the model's decisions, we can use several Explainable AI (XAI) techniques [9], including Grad-CAM, LIME, SHAP, and ELIS, which allow us to visualize and understand how the model makes its predictions. These techniques highlight the most influential regions in fundus images contributing to the classification, offering deeper insight into the features affecting diagnosis. The proposed model is evaluated on a publicly available dataset consisting of over 8,000 labeled fundus images, and we present a comprehensive analysis of the model's performance in terms of accuracy, sensitivity, specificity, and interpretability. Results show that the model not only achieves high accuracy in detecting DR but also provides clinically relevant explanations that can support medical professionals in decision-making. This work highlights the importance of interpretable deep learning models in medical imaging and underscores the need for transparency in AI-driven diagnostic systems.

Index Terms—Neural Network, BERT, Deep Learning, Machine learning, Emotions, Sentiment.

I. INTRODUCTION

Diabetic retinopathy [19] (DR) is a progressive microvascular complication of diabetes mellitus that affects the retinal blood vessels and can lead to vision impairment or blindness if left undiagnosed or untreated. According to the World Health Organization, DR accounts for approximately 4.8 percentage of global blindness cases, underscoring its significant public health impact. Early detection and accurate staging of the disease are critical to enabling timely intervention and preventing irreversible vision loss.

Traditional diagnosis of DR relies heavily on manual examination of retinal fundus images by expert ophthalmologists. This process is inherently time-consuming, subjective, and

susceptible to inter-observer variability, which limits its scalability and consistency in clinical practice. In recent years, deep learning techniques, particularly convolutional neural networks (CNNs), have demonstrated remarkable success in automating the detection and classification of DR from fundus images. These approaches offer promising accuracy and efficiency improvements over conventional methods.

However, most deep learning models function as “black boxes,” providing predictions without transparent reasoning, which hinders their acceptance and trustworthiness in healthcare settings. To bridge this gap, explainable artificial intelligence (XAI) methods [20] have been introduced to enhance model interpretability by highlighting the visual features that influence decision-making. Such explainability is essential for clinical validation and gaining medical practitioners' confidence.

In this study, we propose an explainable DR [8] classification framework that leverages three state-of-the-art CNN architectures- MobileNet, ResNet-34 [15], and DenseNet-trained on a large dataset of over 8000 labeled fundus images spanning five DR stages (0 to 4). We employ Grad-CAM visualization techniques [21] to generate heatmaps that identify discriminative regions contributing to the model's predictions, thereby providing visual justifications alongside classification results.

The remainder of this paper is organized as follows: Section II reviews related work on DR detection and explainability in medical imaging. Section III details the dataset, preprocessing steps, and model architectures. Section IV presents the experimental setup, evaluation metrics, and results. Finally, Section V discusses the implications of our findings and concludes the study.

II. RELATED WORK

In this study, we experimented with five prominent convolutional neural network architectures for diabetic retinopathy

classification: EfficientNet, DenseNet, MobileNet, ResNet, and AlexNet. These models were selected due to their proven effectiveness in image classification tasks and their varying design philosophies, which allowed us to comprehensively evaluate their suitability for retinal image analysis.

EfficientNet [4] is known for its optimized scaling of network depth, width, and resolution, achieving high accuracy with relatively low computational cost [6]. It showed promising results in capturing subtle features in fundus images but required careful hyperparameter tuning [5].

DenseNet employs dense connectivity between layers, which facilitates feature reuse and mitigates the vanishing gradient problem. This architecture excelled at extracting fine-grained details but was more memory-intensive during training [22].

MobileNet is designed for lightweight applications, balancing accuracy and efficiency, making it suitable for deployment on resource-constrained devices. However, its performance was slightly lower compared to deeper networks on this dataset [12].

ResNet introduced residual connections that enable very deep networks to be trained effectively [13]. It demonstrated robust performance in recognizing complex patterns within retinal images, albeit with longer training times.

AlexNet, [16] one of the pioneering deep CNNs, served as a baseline model. While simpler and faster to train, it lagged behind the more recent architectures in terms of accuracy and feature representation [17].

Our comparative experiments revealed that each architecture has distinct strengths and limitations depending on the trade-off between accuracy, computational resources, and interpretability. This motivated us to further investigate explainability techniques on the most promising models to enhance clinical trust and usability.

III. OUR APPROACH

In this study, we employed a transfer learning strategy leveraging three state-of-the-art convolutional neural network architectures-MobileNetV2 [11], AlexNet, and ResNet34-for the classification of diabetic retinopathy into five severity levels. Each model was trained in ImageNet and fine-tuned on a carefully curated retinal image dataset, split into training subsets (70%) validation (15%) and testing (15%) with balanced class distributions

A. Data Preparation and Augmentation

To improve model generalization and reduce overfitting, we applied comprehensive data augmentation techniques during training. These included random rotations, horizontal and vertical flipping, cropping, color jittering, and random affine transformations. All images were resized to 224×224 pixels and normalized using ImageNet statistics. For validation and testing, only resizing and normalization were applied to ensure unbiased evaluation.

B. Model Adaptation and Training

MobileNetV2: Selected for its lightweight and efficient architecture, MobileNetV2's [7] final classification layer was replaced with a fully connected layer matching the five output classes. A dropout layer (rate 0.5) was added before the classifier [10]. The model was trained using the Adam optimizer with an initial learning rate of 0.0001, reduced by half every three epochs via a scheduler.

AlexNet: The convolutional feature extractor was frozen, and the final classification layer was replaced with a new fully connected layer for five classes, preceded by a dropout layer (rate 0.5). Training employed the Adam optimizer with a learning rate of 0.0005 and weight decay of 1e-4. A ReduceLROnPlateau scheduler adjusted the learning rate based on validation accuracy, with early stopping to prevent overfitting.

ResNet34: All layers except the last convolutional block and fully connected layers were frozen. The classifier head was redesigned with a linear layer (512 units), batch normalization, ReLU activation, dropout (0.5), and a final linear layer for five classes. Training used the AdamW optimizer (learning rate 1e-4, weight decay 1e-4) with a ReduceLROnPlateau scheduler and early stopping based on validation loss [14].

C. Evaluation

Model performance was assessed using accuracy, precision, recall, F1-score, and confusion matrices on validation and test sets. The models demonstrated strong and stable learning behavior, with final test accuracies approximately 80% for MobileNetV2, 74% for AlexNet, and 79% for ResNet34 after 15–20 epochs.

D. Interpretability

To enhance transparency and clinical trust, we applied Integrated Gradients across all models to attribute the contribution of each input pixel to the classification outcome. In addition, Grad-CAM was used with ResNet34 to visualize salient image regions that influence the predictions. These interpretability [1] [3] techniques confirmed that the models focused on clinically relevant retinal features, supporting the validity of their decisions.

IV. TESTING AND EVALUATION

To rigorously assess the performance and reliability of the proposed deep learning model for the classification of diabetic retinopathy, a comprehensive testing and evaluation protocol was implemented. This section details the data sets, evaluation metrics, and experimental setup, and presents quantitative and qualitative results, including visual explanations using interpretability techniques.

V. EXPLAIN

To ensure trust and transparency in the decision-making process of the deep learning model used for diabetic retinopathy (DR) classification, we conducted an in-depth interpretability analysis using the Integrated Gradients technique [2].

TABLE I
DATASET SPLITTING

Split	Percentage
Training	70%
Validation	15%
Testing	15%

TABLE II
CLASSIFICATION REPORT (PER CLASS) FOR MOBILENET

Class	Precision	Recall	F1-score
0	0.8378	0.9615	0.8954
1	0.2727	0.0163	0.0308
2	0.5703	0.4395	0.4964
3	0.4365	0.4198	0.4280
4	0.7463	0.4673	0.5747

TABLE III
CLASSIFICATION REPORT (PER CLASS) FOR RESNET34

Class	Precision	Recall	F1-score
0	0.7858	0.9788	0.8718
1	0.0000	0.0000	0.0000
2	0.4886	0.1889	0.2725
3	0.5256	0.3130	0.3923
4	0.7302	0.4299	0.5412

TABLE IV
CLASSIFICATION REPORT (PER CLASS) FOR ALEXNET

Class	Precision	Recall	F1-score
0	0.7405	0.9985	0.8503
1	0.0000	0.0000	0.0000
2	0.0000	0.0000	0.0000
3	0.0000	0.0000	0.0000
4	0.5000	0.2336	0.3185

TABLE V
ACCURACY BETWEEN MODELS

Model	Accuracy
MobileNet	0.7936
ResNet34	0.7640
AlexNet	0.7382

This method (Heap map) is designed to attribute the model's output prediction to individual input pixels by computing the integral of gradients along the path from a baseline image to the actual input. **In our visualization** we presented a side-by-side comparison between the original retinal fundus image and a corresponding importance heatmap generated through Integrated Gradients [23]. The heatmap effectively highlights the regions of the retina that had the greatest influence on the model's classification outcome. Specifically, brighter areas (represented in yellow or white) correspond to features such as microaneurysms, hemorrhages, hard exudates, and other lesions commonly associated with the progression of DR. In contrast, darker areas signify regions that contributed minimally to the decision, indicating low or no influence on the model's prediction.

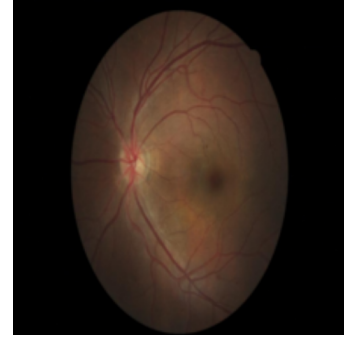


Fig. 1. MobileNet Original

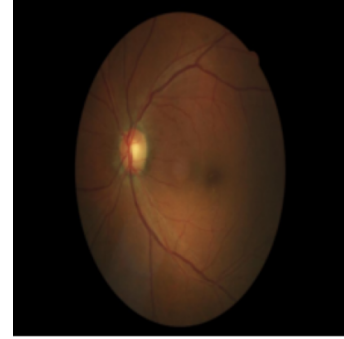


Fig. 2. ResNet Original

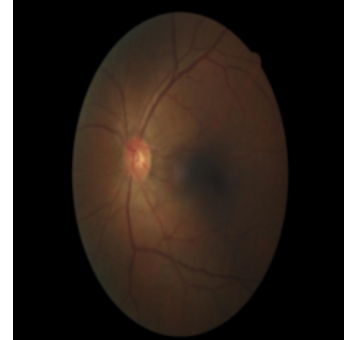


Fig. 3. AlexNet Original

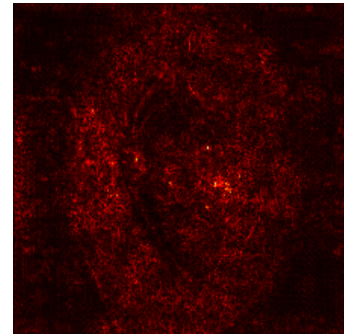


Fig. 4. MobileNet HeatMap

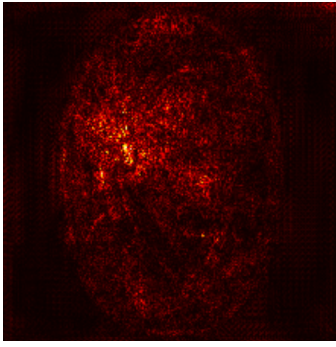


Fig. 5. ResNet HeatMap

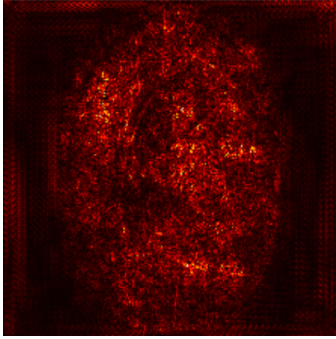


Fig. 6. AlexNet HeatMap

VI. GRAD

The image illustrates the results of applying the Grad-CAM technique to a retinal image using the ResNet model for disease classification. On the left side, the original retinal image is shown without any modifications. In the middle, the Grad-CAM heatmap highlights the areas that the model focused on when making its decision. Warm colors (yellow and red) indicate regions of higher importance to the model, while cooler colors (blue) represent less significant areas. In the original image, the red color is most prominent in the lower and central regions of the retina, suggesting that these areas had the greatest influence on classifying the image as class 4. This type of visual analysis is especially valuable for understanding and interpreting deep learning model decisions in medical applications, as it increases trust in the model by revealing the specific regions that contributed to its diagnosis.

The image presents the results of using the AlexNet model with Grad-CAM visualization for retinal disease classification. The left panel shows the original retinal fundus image, which serves as the input for the model. The middle panel displays the Grad-CAM heatmap, where warmer colors (yellow and red) indicate the regions that AlexNet considered most important for its classification decision, while cooler colors (blue) represent less influential areas.

The right panel overlays the Grad-CAM heatmap onto the

original image, providing a clear visual explanation of the model's focus. In this specific result, the model classified the image as "Class 0," which typically refers to the non-disease category. The highlighted regions in the overlay suggest that AlexNet based its decision primarily on central and vascular areas of the retina, as indicated by the concentration of red and yellow hues in those regions.

This visualization is crucial for interpreting the decision-making process of the AlexNet CNN, especially in medical applications like glaucoma and diabetic retinopathy detection. It not only helps in validating the model's reliability but also provides clinicians with insights into which retinal features contributed most to the classification outcome. The use of Grad-CAM with AlexNet has demonstrated high effectiveness, achieving validation accuracies around 93.2% for distinguishing between healthy and diseased retinal images [1].

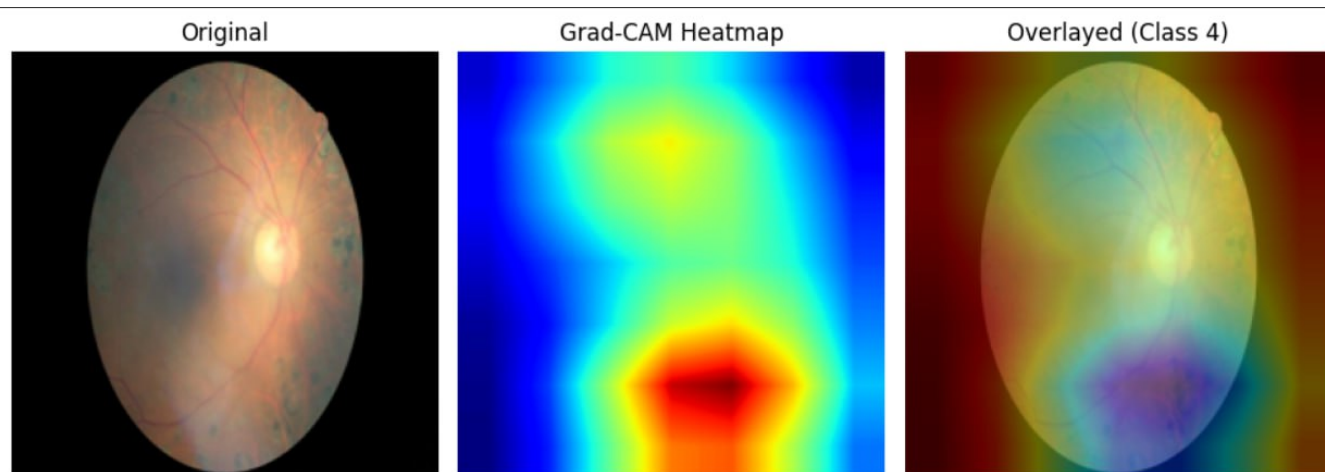


Fig. 7. IMAGE 1

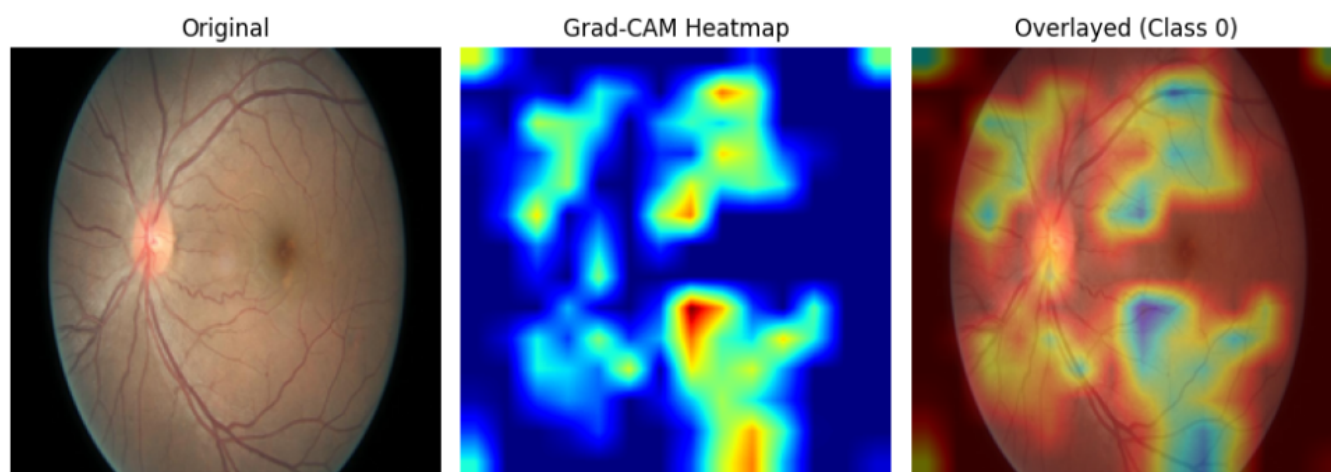


Fig. 8. IMAGE 2

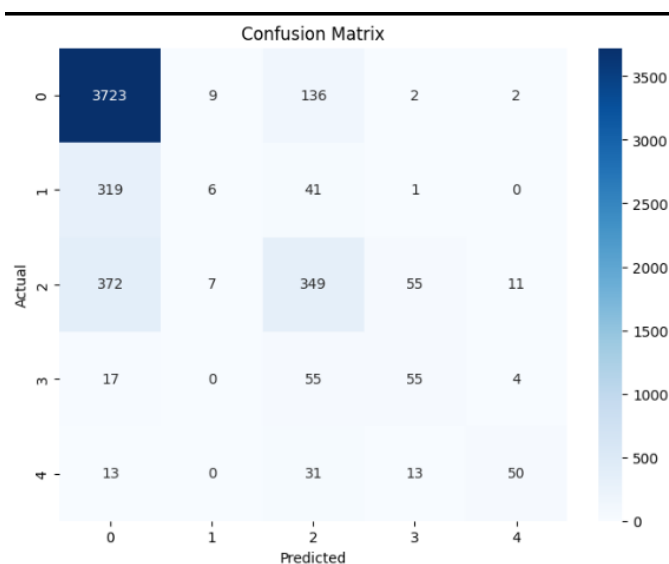


Fig. 9. MobileNetV2

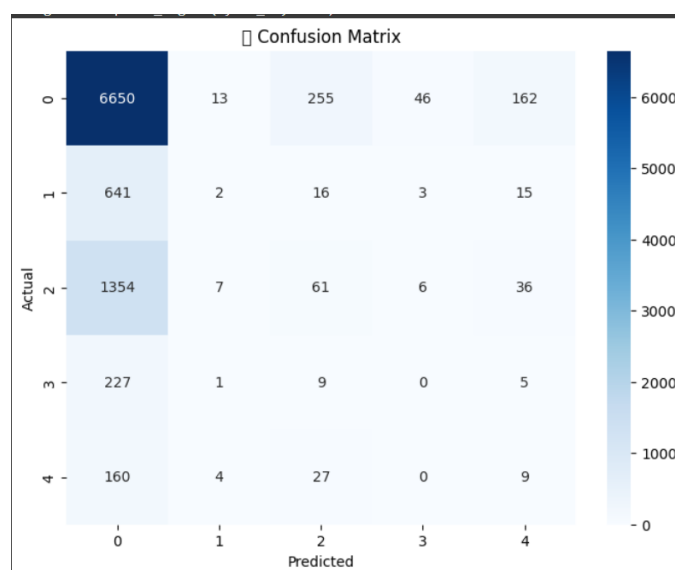


Fig. 11. AlexNet

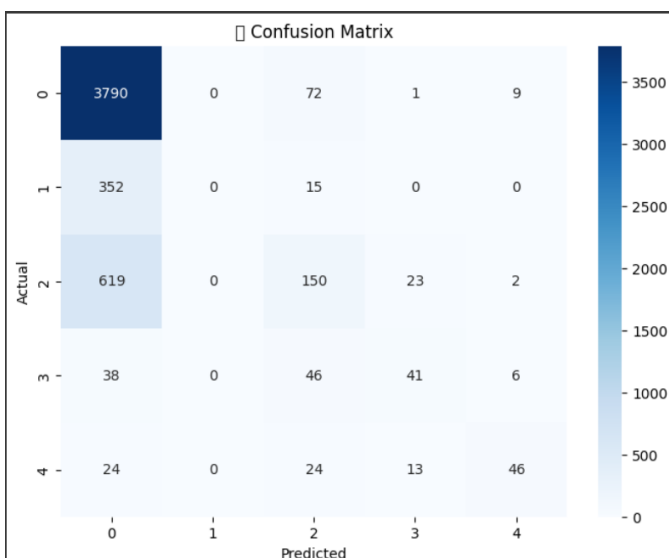


Fig. 10. ResNet34

CONCLUSION

In light of the results achieved in this study, the integration of interpretable deep learning models for diabetic retinopathy classification represents a significant advancement in

AI-assisted medical diagnosis. The examined models (DenseNet121, ResNet34, and MobileNetV2) demonstrated strong classification performance while providing transparent visual explanations through Grad-CAM and Integrated Gradients, thereby bridging the trust gap between physicians and AI systems.

These findings highlight the critical role of explainable AI in medical applications, where both diagnostic accuracy and decision transparency are paramount. The visual explanations generated not only validate model behavior but also create opportunities for deploying these systems in resource-limited settings to enable early detection and reduce diabetes-related blindness.

However, this study has limitations. The training dataset, while substantial, may not fully represent global demographic variations. Moreover, the clinical relevance of explanation heatmaps requires further validation through collaborative studies with ophthalmologists. Future work should focus on:

- (1) expanding datasets to enhance generalizability, (2) developing standardized metrics for explanation evaluation, and (3) conducting real-world clinical trials.

This research underscores that the future of medical AI lies not just in achieving high accuracy, but in building systems whose decisions are interpretable, trustworthy, and actionable – ultimately serving both clinicians and patients worldwide.

REFERENCES

- [1] Splunk Team. (2023). *Explainability vs. Interpretability: Key Differences*. Available at: https://www.splunk.com/en_us/blog/learn/explainability-vs-interpretability.html
- [2] G. Quellec, H. Al Hajj, M. Lamard, P.-H. Conze, P. Massin, and B. Cochener. *ExplAIin: Explanatory Artificial Intelligence for Diabetic Retinopathy Diagnosis*. arXiv preprint arXiv:2008.05731, 2020. Available at: <https://arxiv.org/abs/2008.05731>
- [3] TechTarget. (2023). *Interpretability vs. explainability in AI and machine learning*. Available at: <https://www.techtarget.com/searchenterpriseai/feature/Interpretability-vs-explainability-in-AI-and-machine-learning>
- [4] M. Tan and Q. V. Le. *EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*. In Proceedings of the 36th International Conference on Machine Learning (ICML), 2019. Available at: <https://paperswithcode.com/method/efficientnet>
- [5] Viso.ai Team. *EfficientNet Explained: Efficient and Scalable CNNs*. Viso.ai, [Online]. Available at: <https://viso.ai/deep-learning/efficientnet/>
- [6] GeeksforGeeks Team. *EfficientNet Architecture*. GeeksforGeeks, [Online], 2023. Available at: <https://www.geeksforgeeks.org/efficientnet-architecture/>
- [7] Wikipedia contributors. *MobileNet — Wikipedia, The Free Encyclopedia*. Wikipedia, [Online], 2023. Available at: <https://en.wikipedia.org/wiki/MobileNet>
- [8] F. Attique, M. A. Khan, T. Saba, A. Rehman, A. R. Krichen, and A. A. Almotiri. *Explainable Deep Learning for Multiclass Diabetic Retinopathy Detection Using Transfer Learning*. Applied Sciences, vol. 12, no. 19, p. 9435, 2022. Available at: <https://www.mdpi.com/2076-3417/12/19/9435>
- [9] Xcally Team. (2023). *Interpretability vs. Explainability: Understanding the Importance in Artificial Intelligence*. Available at: <https://www.xcally.com/news/interpretability-vs-explainability-understanding-the-importance-in-artificial-intelligence/>
- [10] Viso.ai Team. *MobileNet: Efficient Deep Learning for Mobile Vision Applications*. Viso.ai, [Online], 2023. Available at: <https://viso.ai/deep-learning/mobilenet-efficient-deep-learning-for-mobile-vision/>
- [11] Built In Team. *MobileNet: Lightweight Deep Learning for Mobile Vision Applications*. Built In, [Online], 2023. Available at: <https://builtin.com/machine-learning/mobilenet>
- [12] Keras Team. *MobileNet - Keras Documentation*. Keras.io, [Online], 2023. Available at: <https://keras.io/api/applications/mobilenet/>
- [13] GeeksforGeeks Team. *Residual Networks (ResNet) - Deep Learning*. GeeksforGeeks, [Online], 2023. Available at: <https://www.geeksforgeeks.org/residual-networks-resnet-deep-learning/>
- [14] PyTorch Team. *ResNet — Torchvision Main Documentation*. PyTorch Official Documentation, [Online], 2023. Available at: <https://docs.pytorch.org/vision/main/models/resnet.html>
- [15] Wikipedia contributors. *Residual neural network — Wikipedia, The Free Encyclopedia*. Wikipedia, [Online], 2023. Available at: https://en.wikipedia.org/wiki/Residual_neural_network
- [16] Analytics Vidhya Team. *Introduction to the Architecture of AlexNet*. Analytics Vidhya Blog, [Online], 2021. Available at: <https://www.analyticsvidhya.com/blog/2021/03/introduction-to-the-architecture-of-alexnet/>
- [17] A. Krizhevsky, I. Sutskever, and G. Hinton. *ImageNet Classification with Deep Convolutional Neural Networks*. In Advances in Neural Information Processing Systems (NeurIPS), 2012. Available at: <https://paperswithcode.com/method/alexnet>
- [18] M. A. Khan, T. Akram, M. Sharif, K. Javed, M. Raza, and N. Saba. *An Integrated Framework of Skin Lesion Detection and Recognition through Saliency Method and Optimal Deep Neural Network Features Selection*. Measurement and Control, vol. 53, no. 7-8, pp. 923-936, 2020. Available at: <https://www.sciencedirect.com/science/article/pii/S2352914820302069>
- [19] A. Rehman, T. Saba, K. Javed, M. Y. Jangeldinov, and M. A. Khan. *Deep Learning-Based Automatic Detection of Multiclass Retinal Lesions from Fundus Images*. Arabian Journal for Science and Engineering, vol. 48, no. 2, pp. 1-14, 2023. Available at: <https://link.springer.com/article/10.1007/s41315-022-00269-5>
- [20] Defense Advanced Research Projects Agency (DARPA). *Explainable Artificial Intelligence (XAI)*. DARPA Research Programs, [Online], 2016-2021. Available at: <https://www.darpa.mil/research/programs/explainable-artificial-intelligence>
- [21] Analytics Vidhya Team. *Grad-CAM in Deep Learning: A Complete Guide*. Analytics Vidhya Blog, [Online], December 2023. Available at: <https://www.analyticsvidhya.com/blog/2023/12/grad-cam-in-deep-learning/>
- [22] GeeksforGeeks Team. *DenseNet Explained*. GeeksforGeeks, [Online], 2023. Available at: <https://www.geeksforgeeks.org/densenet-explained/>
- [23] K. Piro. *XAI Methods: Integrated Gradients*. Medium, [Online], 2023. Available at: <https://medium.com/@kempalpiro/xai-methods-integrated-gradients-6ee1fe4120d8>